



情報知識ネットワーク特論 「情報検索とパターン照合」

情報科学研究科 コンピュータサイエンス専攻
情報知識ネットワーク研究室

喜田拓也

第5回

正規表現の照合

正規表現について
照合処理のながれ
構文木(parse tree)の構築
NFAの構築
NFAのシミュレーション手法



正規表現とは？

■ 柔軟で強力なテキスト照合のための記法

– ファイル名の正規表現の例:

> rm *.txt

“任意のファイル名”.txt にマッチ

> cp Important[0-9].doc

Important0.doc～Important9.doc にマッチ

– 検索ツールGrepの正規表現の例:

> grep -E “for.+(256|CHAR_SIZE)” *.c

– プログラミング言語Perlの正規表現の例:

\$line = m|^http://.+¥jp/.+|

“http://”で始まり、“.jp/”を含む
文字列にマッチ

■ 正規表現は任意の正規集合(正規言語)を表現できる

– 有限オートマトンが受理できる言語(文字列の集合)Lを表現



正規表現の定義

■ 定義:

正規表現とは、 $\Sigma \cup \{\varepsilon, |, \cdot, *, (,)\}$ 上の文字列であり、以下の規則で再帰的に定義される。

- (1) ε と Σ の要素は正規表現である
- (2) α と β が正規表現ならば $(\alpha \cdot \beta)$ も正規表現である
- (3) α と β が正規表現ならば $(\alpha | \beta)$ も正規表現である
- (4) α が正規表現ならば α^* も正規表現である
- (5) 上から導かれるものだけが正規表現である

例: $(A \cdot ((A \cdot T) | (C \cdot G))^*) \rightarrow A(AT|CG)^*$

簡略のために、 $(\alpha \cdot \beta)$ を単に $\alpha \beta$ と記述する

※ ‘|’, ‘.’, ‘*’ は、オペレータ (operator) と呼ばれる
また、‘+’ は、 α を正規表現とすると、 $\alpha^+ = \alpha \cdot \alpha^*$ の意味で使用されることがある



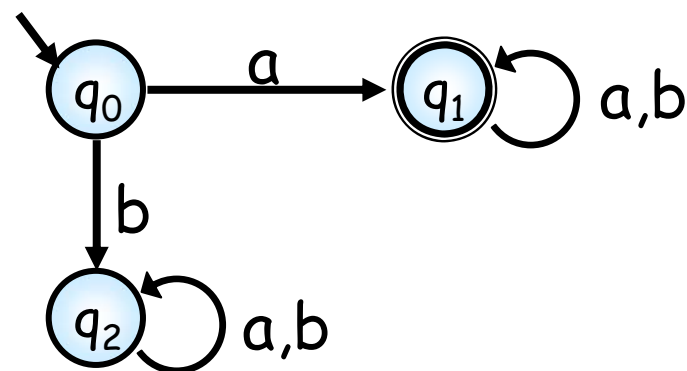
正規表現の意味づけ

- 正規表現を Σ^* の部分集合 (言語L) に写像する
 - (i) $\|\varepsilon\| = \{\varepsilon\}$
 - (ii) $a \in \Sigma$ に対して $\|a\| = \{a\}$
 - (iii) 正規表現 α, β に対して $\|(\alpha \cdot \beta)\| = \|\alpha\| \cdot \|\beta\|$
 - (iv) 正規表現 α, β に対して $\|(\alpha | \beta)\| = \|\alpha\| \cup \|\beta\|$
 - (v) 正規表現 α に対して $\|\alpha^*\| = \|\alpha\|^*$

- 例: $(a \cdot (a | b)^*)$

$$\begin{aligned}
 & \| (a \cdot (a | b)^*) \| \\
 &= \|a\| \cdot \|(a | b)^*\| \\
 &= \{a\} \cdot \|(a | b)\|^* \\
 &= \{a\} \cdot (\{a\} \cup \{b\})^* \\
 &= \{ax \mid x \in \{a, b\}^*\}
 \end{aligned}$$

左の例と等価な決定性オートマトン (DFA)

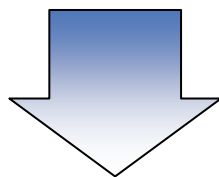


演習: $(AT|GA)(TT)^*$ はどのような言語になるだろうか?



正規表現を照合するとは？

- 正規表現の照合問題：
 - 正規表現 α で定義される言語 $L(\alpha) = \|\alpha\|$ に含まれる任意の文字列をテキスト中から探し出す問題
- 正規表現と有限オートマトンは言語を定義する能力が等しい！
 - 正規表現で表現できる言語を受理する有限オートマトンを構築できる
 - 逆に、有限オートマトンが受理する言語を表現する正規表現も存在する
 - ※ 有川節夫・宮野悟著、「オートマトンと計算可能性」参照(2.5 正規表現と正規集合)

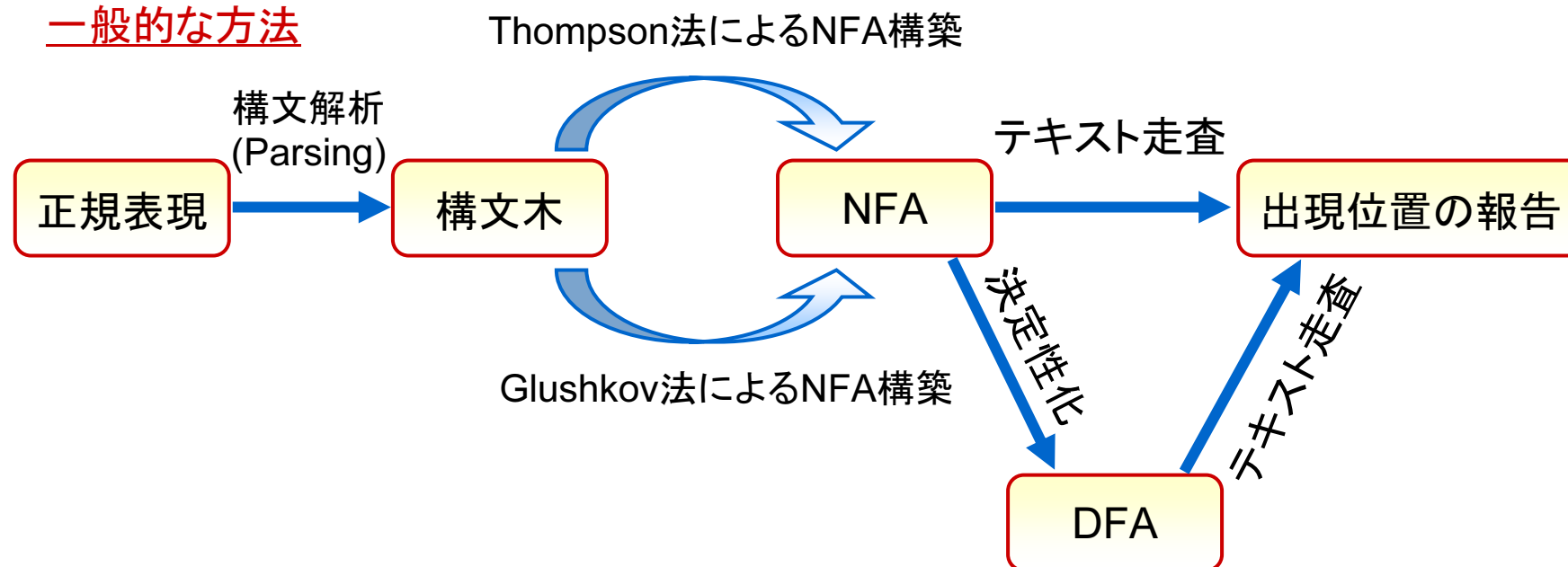


- 正規表現に対応する(非)決定性オートマトンを作成し、その動きをシミュレートすればよい。
 - DFAよりも、NFAへ変換するほうが容易
 - ただし、初期状態は常にアクティブ
 - テキストを読んでいく過程で、オートマトンが受理状態に到達したらパタン出現



照合処理のながれ

一般的な方法



Filter手法による方法

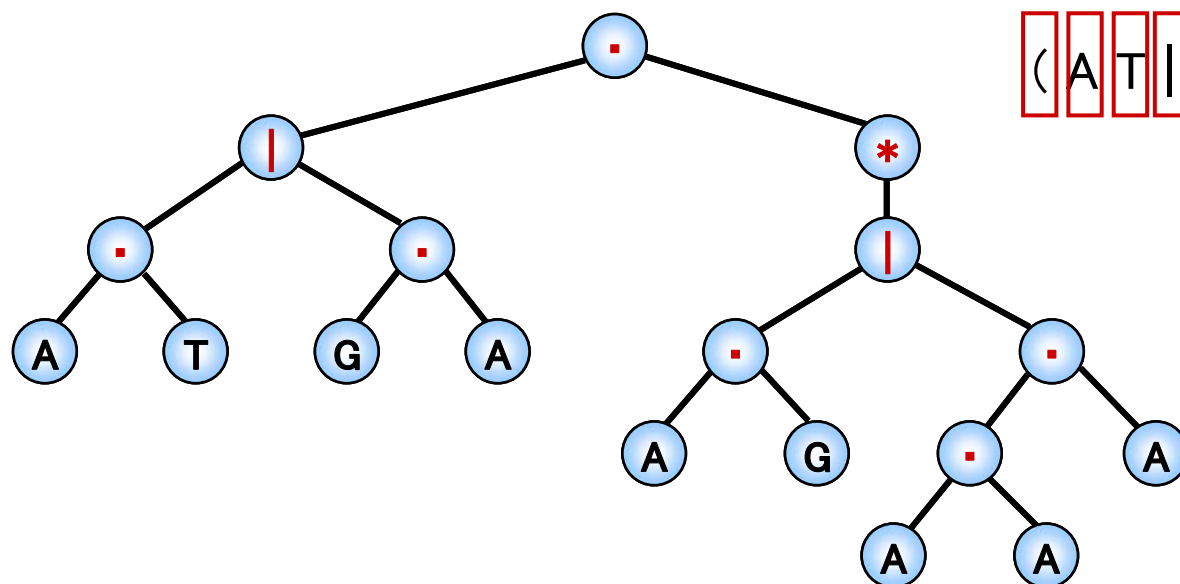




構文木 (parse tree) の構築

- 構文木: NFAを作りやすくするための準備として用いる木構造
 - 葉ノードは各々アルファベット Σ 上の文字 $a \in \Sigma$ もしくは空語 ε でラベル付けされる
 - 内部ノードはオペレータ $\{ |, \cdot, * \}$ でラベル付けされる
 - LexやFlexなどの構文解析ツールでも変換できるが、正規表現のパーサとしては大げさすぎる (→ 次のスライドの擬似コードで十分)

例: 正規表現 $RE=(AT|GA)((AG|AAA)^*)$ の構文木 T_{RE}



$(AT|GA)((AG|AAA)^*)$

括弧の深さ	オペレータ
1	
2	



擬似コード

```

Parse (p=p1p2...pm, last)
1  v ← θ;
2  while plast ≠ $ do
3      if plast ∈ Σ or plast = ε then /* normal character */
4          vr ← Create a node with plast;
5          if v ≠ θ then v ← [·](v, vr);
6          else v ← vr;
7          last ← last + 1;
8      else if plast = '|' then /* union operator */
9          (vr, last) ← Parse(p, last + 1);
10         v ← [|](v, vr);
11     else if plast = '*' then /* star operator */
12         v ← [*](v);
13         last ← last + 1;
14     else if plast = '(' then /* open parenthesis */
15         (vr, last) ← Parse(p, last + 1);
16         last ← last + 1;
17         if v ≠ θ then v ← [·](v, vr);
18         else v ← vr;
19     else if plast = ')' then /* close parenthesis */
20         return (v, last);
21     end of if
22 end of while
23 return (v, last);

```



ThompsonのNFA構築法

K. Thompson. Regular expression search algorithm. *Communications of the ACM*, 11:419-422, 1968.

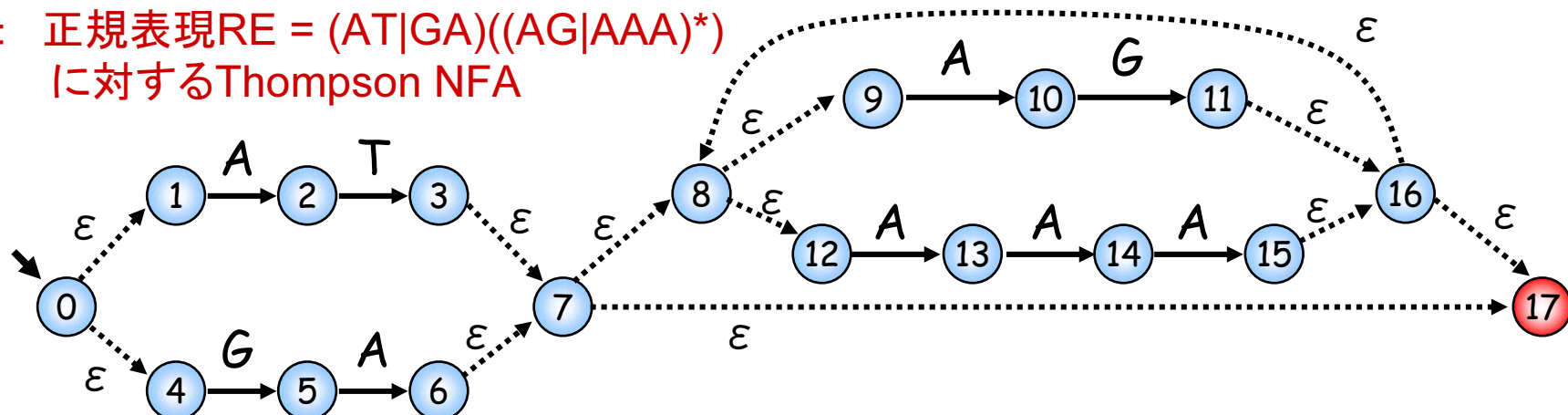
■ アイデア

- 正規表現RE の構文木 T_{RE} をpost-orderで巡回しながら、各ノード v を頂点とする部分木に対応する言語 $L(RE_v)$ を受け取るオートマトン $Th(v)$ を構築していく
- $Th(v)$ は、 v の子供を頂点とする部分木に対するオートマトンどうしを、 ϵ 遷移で連結することで得られるのがポイント

■ Thompson NFAの性質

- 状態数 $< 2m$ 、状態遷移の数 $< 4m \rightarrow O(m)$
- ϵ 遷移を含む
- ϵ 遷移以外の遷移は必ず i 番目から $i+1$ 番目に遷移する

例: 正規表現RE = (AT|GA)((AG|AAA)*)
に対するThompson NFA

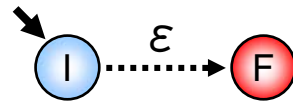




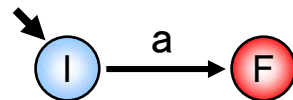
NFA構築アルゴリズム

- 構文木 TREに対して、post-orderでノードを探索しつつ、各ノードについて次のようにオートマトンを生成・連結する

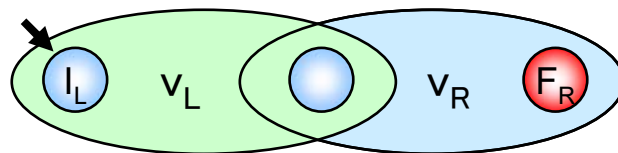
(i) ノード v が空語 ε の場合



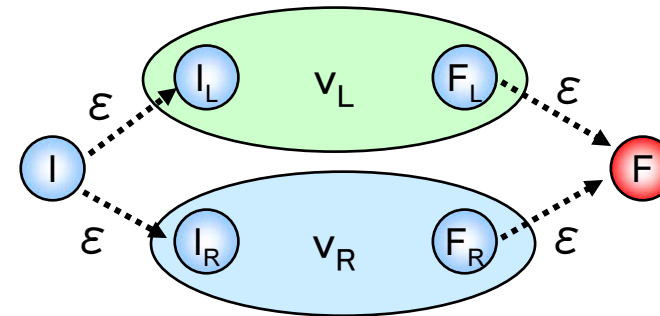
(ii) ノード v が文字 a の場合



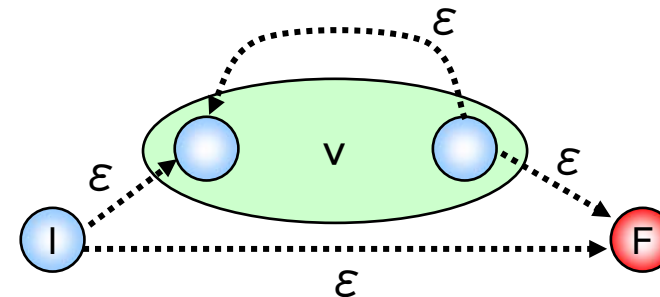
(iii) ノード v が連結 \cdot の場合
 $\rightarrow (v_L \cdot v_R)$



(iv) ノード v が選択 $|$ の場合 $\rightarrow (v_L | v_R)$



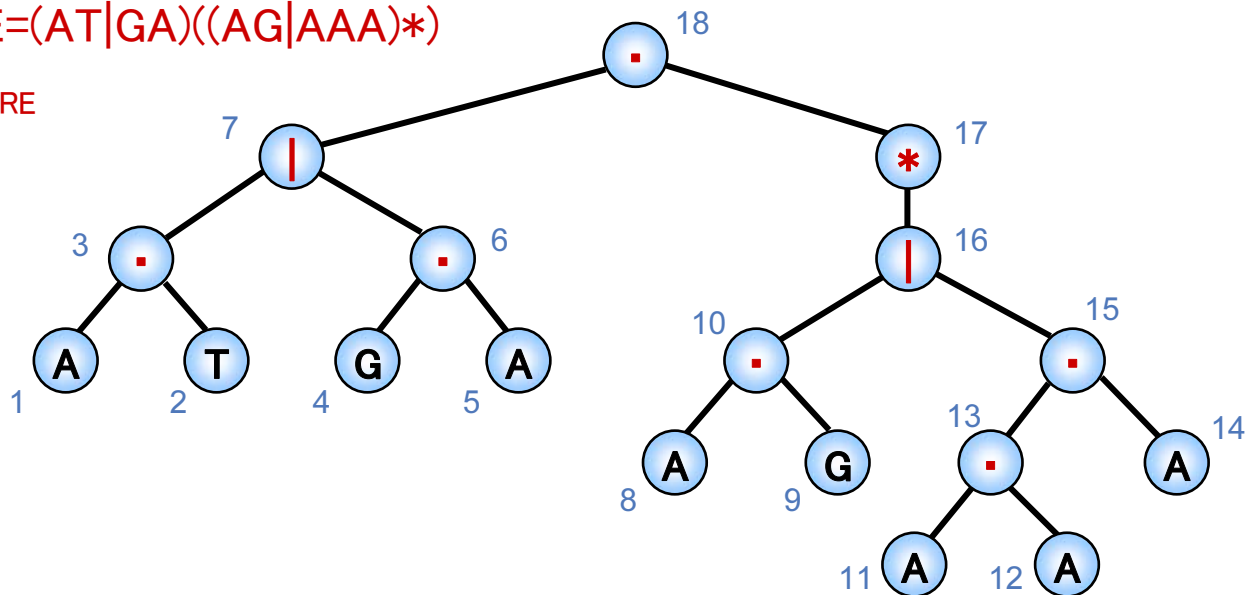
(v) ノード v が繰り返し $*$ の場合 $\rightarrow v^*$



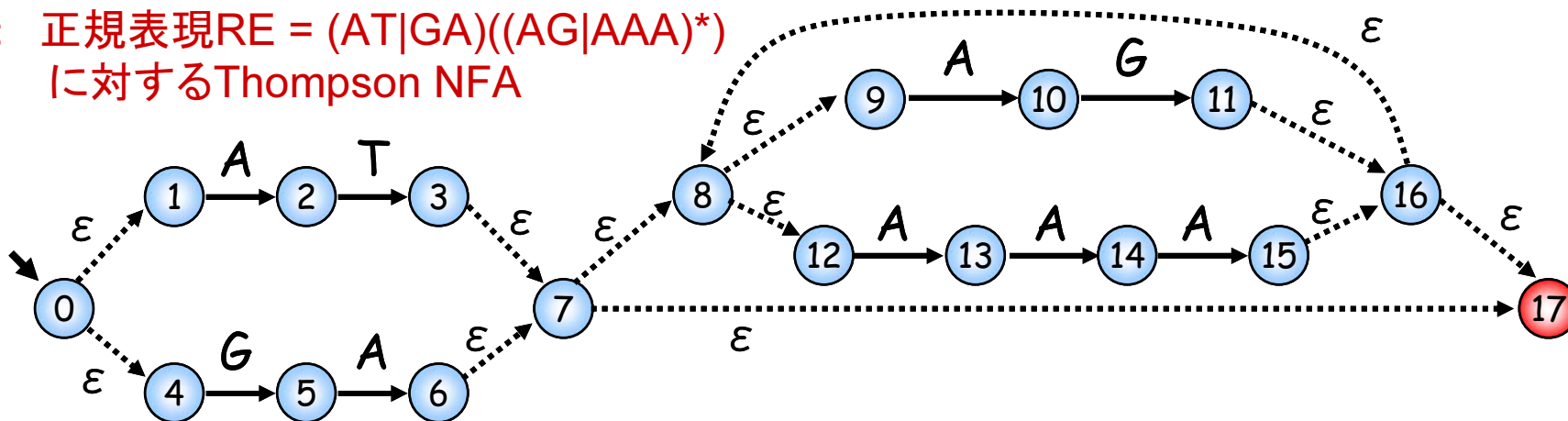


NFA構築アルゴリズムの動作

例：正規表現 $RE=(AT|GA)((AG|AAA)^*)$
の構文木 T_{RE}



例：正規表現 $RE = (AT|GA)((AG|AAA)^*)$
に対するThompson NFA





擬似コード

Thompson_recur (v)

```
1  if v = "|" (vL, vR) or v = "." (vL, vR) then
2      Th(vL) ← Thompson_recur(vL);
3      Th(vR) ← Thompson_recur(vR);
4  else if v = "*" (vC) then Th(v) ← Thompson_recur(vC);
5  /* ここまでが再帰的な処理 (post order traverse) */
6  if v = (ε) then return construction (i);
7  if v = (α), α ∈ Σ then return construction (ii);
8  if v = "." (vL, vR) then return construction (iii);
9  if v = "|" (vL, vR) then return construction (iv);
10 if v = "*" (vC) then return construction (v);
```

Thompon(RE)

```
11 vRE ← Parse(RE$, 1); /* 構文木を構築する */
12 Th(vRE) ← Thompson_recur(vRE);
```



GlushkovのNFA構築法

V-M. Glushkov. The abstract theory of automata. *Russian Mathematical Surveys*, 16:1-53, 1961.

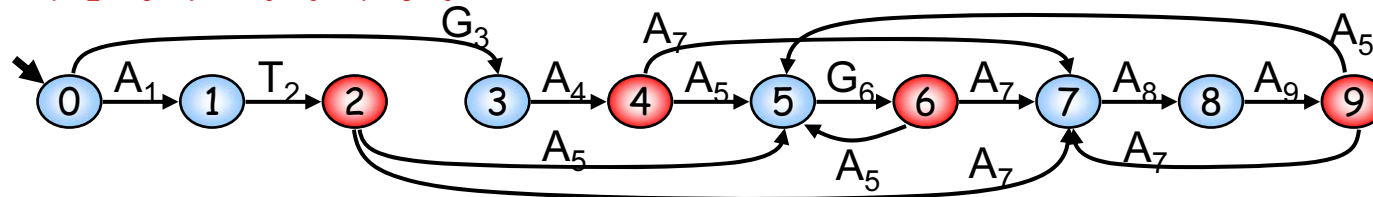
■ アイデア

- 正規表現RE中の各文字 $a \in \Sigma$ に前から順に番号を振り、これをRE'とする
(番号付きのアルファベットを Σ' とする)
 - 例: $RE = (AT|GA)((AG|AAA)^*) \rightarrow RE' = (A_1T_2|G_3A_4)((A_5G_6|A_7A_8A_9)^*)$
- 言語 $L(RE')$ を受理するNFAを作り、番号を取り除いて最終的なNFAを得る

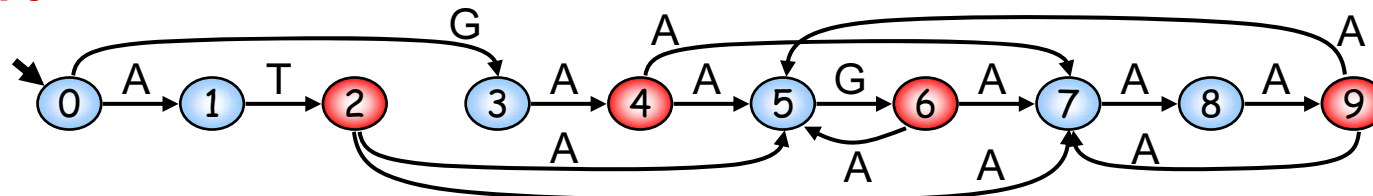
■ Glushkov NFAの性質

- 状態数はちょうど $m+1$ 個、状態遷移の数は $O(m^2)$
- ε 遷移を含まない
- 任意のノードについて、そのノードに入ってくる遷移のラベルはすべて等しい

例: $RE' = (A_1T_2|G_3A_4)((A_5G_6|A_7A_8A_9)^*)$ に対するNFA



例: 最終的なGlushkov NFA





NFA構築アルゴリズム(1)

■ 構築手順:

- 正規表現RE中の各文字 $a \in \Sigma$ に前から順に番号を振り、これをRE'とする
 - $\text{Pos}(RE') = \{1 \cdots m\}$ 、 Σ' : 番号付けされたアルファベット
- 構文木 $T_{RE'}$ をpost-orderで巡回しながら、各ノード v を頂点とする部分木に対応する言語 RE'_v について集合 $\text{First}(RE'_v)$ 、 $\text{Last}(RE'_v)$ 、関数 Empty_v 、および、RE'の位置 x についての関数 $\text{Follow}(RE', x)$ を計算する

- $\text{First}(RE') = \{x \in \text{Pos}(RE') \mid \exists u \in \Sigma'^*, \alpha_x u \in L(RE')\}$

NFAが開始される位置

- $\text{Last}(RE') = \{x \in \text{Pos}(RE') \mid \exists u \in \Sigma'^*, u \alpha_x \in L(RE')\}$

NFAの最終状態の位置

- $\text{Follow}(RE', x) = \{y \in \text{Pos}(RE') \mid \exists u, v \in \Sigma'^*, u \alpha_x \alpha_y v \in L(RE')\}$

遷移関数のつながり

- Empty_{RE} : ε が $L(RE)$ に属するなら $\{\varepsilon\}$ 、そうでないならば ϕ を返す関数

これは、次のようにして再帰的に計算できる

$$\text{Empty}_\varepsilon = \{\varepsilon\},$$

$$\text{Empty}_{a \in \Sigma} = \phi,$$

$$\text{Empty}_{RE1|RE2} = \text{Empty}_{RE1} \cup \text{Empty}_{RE2},$$

$$\text{Empty}_{RE1 \cdot RE2} = \text{Empty}_{RE1} \cap \text{Empty}_{RE2},$$

$$\text{Empty}_{RE^*} = \{\varepsilon\}.$$

NFAの初期状態が
終状態かどうか

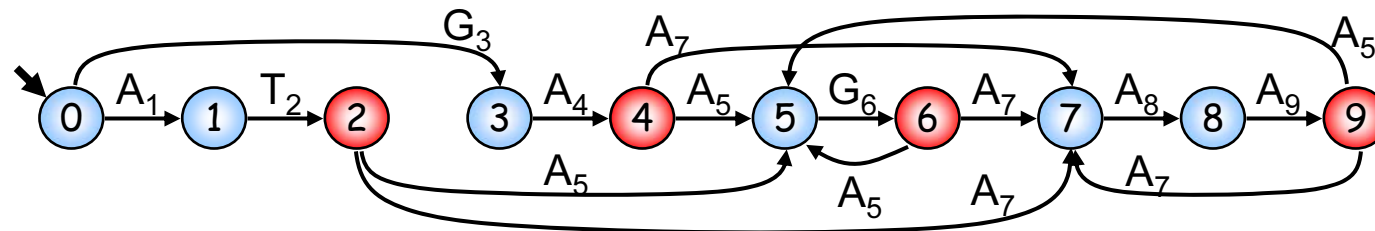
- 上で得られた値をもとに、NFAを構築する



NFA構築アルゴリズム(2)

- 言語 $L(RE')$ を受理する (Glushkov) NFA $GL' = (S, \Sigma', I, F, \delta')$
 - S : 状態の集合. $S = \{0, 1, \dots, m\}$
 - Σ' : 番号付けされた Σ
 - I : 初期状態、 $I = 0$
 - F : 最終状態
 $F = \text{Last}(RE') \cup (\text{Empty}_{RE} \cdot \{0\})$.
 - δ' : 以下で定義される遷移関数
 $\forall x \in \text{Pos}(RE'), \forall y \in \text{Follow}(RE', x), \delta'(x, \alpha_y) = y$
 初期状態からの遷移は次のとおり
 $\forall y \in \text{First}(RE'), \delta'(0, \alpha_y) = y$

例: $RE' = (A_1 T_2 | G_3 A_4) ((A_5 G_6 | A_7 A_8 A_9)^*)$ に対する NFA





擬似コード

```

Glushkov_variables ( $v_{RE}$ , lpos)
1  if  $v = [\bar{\quad}](v_l, v_r)$  or  $v = [\cdot](v_l, v_r)$  then
2      lpos  $\leftarrow$  Glushkov_variables( $v_l$ , lpos);
3      lpos  $\leftarrow$  Glushkov_variables( $v_r$ , lpos);
4  else if  $v = [*](v_*)$  then lpos  $\leftarrow$  Glushkov_variables( $v_*$ , lpos);
5  end of if
6  if  $v = (\varepsilon)$  then
7      First( $v$ )  $\leftarrow \phi$ , Last( $v$ )  $\leftarrow \phi$ , Empty $_v$   $\leftarrow \{\varepsilon\}$ ;
8  else if  $v = (a)$ ,  $a \in \Sigma$  then
9      lpos  $\leftarrow$  lpos + 1;
10     First( $v$ )  $\leftarrow \{lpos\}$ , Last( $v$ )  $\leftarrow \{lpos\}$ , Empty $_v$   $\leftarrow \phi$ , Follow(lpos)  $\leftarrow \phi$ ;
11  else if  $v = [ | ](v_l, v_r)$  then
12     First( $v$ )  $\leftarrow$  First( $v_l$ )  $\cup$  First( $v_r$ );
13     Last( $v$ )  $\leftarrow$  Last( $v_l$ )  $\cup$  Last( $v_r$ );
14     Empty $_v$   $\leftarrow$  Empty $_{v_l}$   $\cup$  Empty $_{v_r}$ ;
15  else if  $v = [\cdot](v_l, v_r)$  then
16     First( $v$ )  $\leftarrow$  First( $v_l$ )  $\cup$  (Empty $_{v_l} \cdot$  First( $v_r$ ));
17     Last( $v$ )  $\leftarrow$  (Empty $_{v_r} \cdot$  Last( $v_l$ ))  $\cup$  Last( $v_r$ );
18     Empty $_v$   $\leftarrow$  Empty $_{v_l} \cap$  Empty $_{v_r}$ ;
19     for  $x \in$  Last( $v_l$ ) do Follow( $x$ )  $\leftarrow$  Follow( $x$ )  $\cup$  First( $v_r$ );
20  else if  $v = [*](v_*)$  then
21     First( $v$ )  $\leftarrow$  First( $v_*$ ), Last( $v$ )  $\leftarrow$  Last( $v_*$ ), Empty $_v$   $\leftarrow \{\varepsilon\}$ ;
22     for  $x \in$  Last( $v_*$ ) do Follow( $x$ )  $\leftarrow$  Follow( $x$ )  $\cup$  First( $v_*$ );
23  end of if
24  return lpos;

```

トータル $O(m^3)$ 時間

$O(m^2)$ 時間かかる



擬似コード(続き)

Glushkov (RE)

```
1  /* 正規表現をパースして構文木を作る */
2  vRE ← Parse(RE$, 1);
3
4  /* 構文木を使って各変数を計算する */
5  m ← Glushkov_variables(vRE, 0);
6
7  /* 計算した変数を使ってNFA GL(S, Σ, I, F, δ)を構築する */
8  Δ ← φ;
9  for i ∈ 0...m do create state I;
10 for x ∈ First(vRE) do Δ ← Δ ∪ {(0, αx, x)};
11 for i ∈ 0...m do
12     for i ∈ Follow(i) do Δ ← Δ ∪ {(i, αx, x)};
13 end of for
14 for x ∈ Last(vRE) ∪ (EmptyvRE · {0}) do mark x as terminal;
```



ちょっと、ひといき・・・

- ここまでのまとめ
 - 正規表現：有限オートマトンと言語を定義する能力が等しい
 - 正規表現の照合処理のながれ
 - 構文木 (parse tree) を構築してからNFAへ変換
 - NFAをシミュレートしてテキストを走査
 - NFAの構築方法
 - ThompsonのNFA → 状態数 $< 2m$ 、状態遷移数 $< 4m$
 - GlushkovのNFA → 状態数 $= m+1$ 、状態遷移数 $O(m^2)$



フンボルトペンギン (旭山動物園にて '05.3.13)

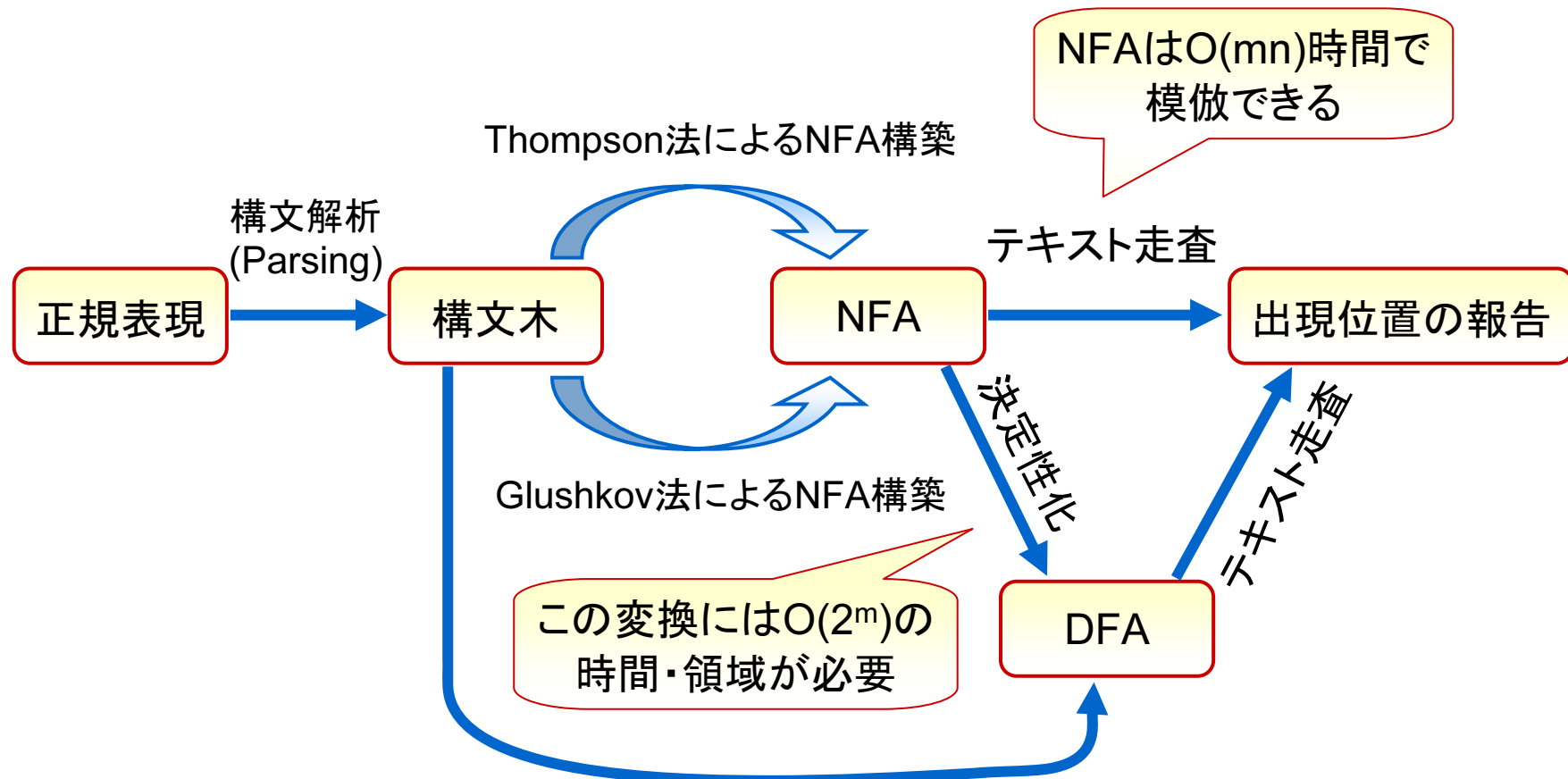
～トリビア～

POSIX準拠の正規表現は正規表現？

現在、正規表現を処理するプログラムの多くはNFAの動作を「バックトラック」を使って実装している。「後方参照」は、大抵の場合この技法の上で実現されるが、明らかに「正規表現」の枠組みからは逸脱している。→ 例: $(a+)b\yen1$



照合処理のながれ(再)



実はDFAへ直接に変換する方法もある

※A. V. Aho, R. Sethi, and J. D. Ullman. Compilers – Principles, Techniques and Tools. Addison-Wesley, 1986. (邦訳: コンパイラ 原理・技法・ツール)を参照 (3.9節)



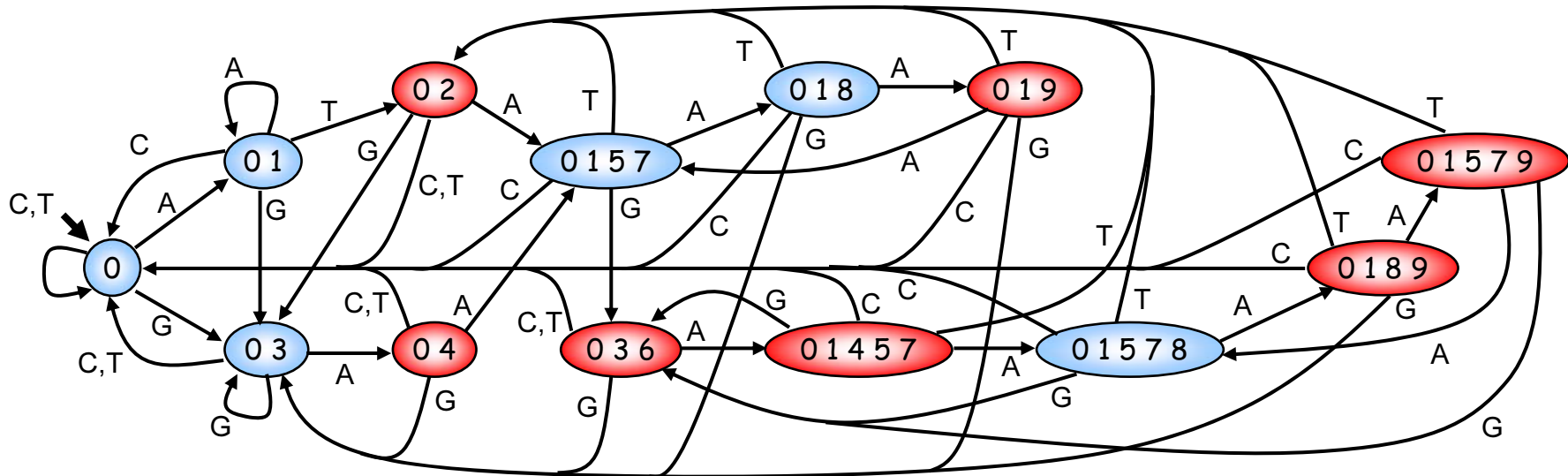
NFAのシミュレーション手法

- Thompson提案のNFAシミュレーション手法
 - 最も簡単な手法
 - $O(m)$ の大きさのリストでactiveな状態を保持し、文字ごとのNFAの状態更新を $O(m)$ で行う。
 - 明らかに $O(mn)$ 時間かかる
- 等価なDFAに変換してシミュレートする
 - 古典的に用いられている手法
 - A. V. Aho, R. Sethi, and J. D. Ullman. *Compilers - Principles, Techniques and Tools*. Addison-Wesley, 1986. (邦訳:コンパイラ 原理・技法・ツール)を参照
 - 前処理として変換する $\rightarrow O(2^m)$ 時間・領域かかる
 - テキストを走査する際、動的に構築する手法もある
- ハイブリッドな手法
 - E. W. Myers. A four russians algorithm for regular expression pattern matching. *Journal of the ACM*, 39(2):430-448, 1992.
 - NFAとDFAを組み合わせて効率をあげる手法
 - ThompsonのNFAを $O(k)$ 個のノードのモジュールに分割し、モジュール毎にDFA化する。各モジュール間の遷移はNFAとしてシミュレートする。
- Bit-parallel手法による高速なNFAシミュレーション
 - ThompsonのNFAをシミュレート: S. Wu and U. Manber[1992]の手法
 - GlushkovのNFAをシミュレート: G. Navarro and M. Raffinot[1999]の手法、ほか



等価なDFAに変換してシミュレートする

例: RE = (AT|GA)((AG|AAA)*)に対するGlushkov NFAから変換したDFA



DFA Classical ($N = (Q, \Sigma, l, F, \Delta)$, $T = t_1 t_2 \dots t_n$)

1 Preprocessing:

2 **for** $\sigma \in \Sigma$ **do** $\Delta \leftarrow \Delta \cup (i, \sigma, l)$;

3 $(Q_d, \Sigma, l_d, F_d, \delta) \leftarrow \mathbf{BuildDFA}(N)$; /* NFA Nと等価なDFAを作成 */

4 Searching:

5 $s \leftarrow l_d$;

6 **for** $pos \in 1 \dots n$ **do**

7 **if** $s \in F_d$ **then** report an occurrence ending at $pos - 1$;

8 $s \leftarrow \delta(s, t_{pos})$;

9 **end of for**



Bit-parallel Thompson

S. Wu and U. Manber. Fast text searching allowing errors. *Communications of the ACM*, 35(10):83-91, 1992.

- Thompson NFAをbit-parallelでシミュレートする
 - Thompson NFAでは、(ϵ 遷移を除くと) i 番目の状態の次は $i+1$ 番目
→ Shift-And法に似たbit-parallel化ができる!
 - ϵ 遷移については、別途シミュレートする
 - 大きさ 2^L のマスクテーブルが必要 (L はNFAの状態数)
 - 前処理全体で $O(2^L + m|\Sigma|)$ 時間かかる
 - L が十分小さいときは $O(n)$ 時間でテキストを走査できる
- NFA $G_L = (Q = \{s_0, \dots, s_{|Q|-1}\}, \Sigma, I = s_0, F, \Delta)$ について
 - NFAのマスクビット表現: $Q_n = \{0, \dots, |Q|-1\}$, $I_n = 0^{|Q|-1}1$, $F_n = \bigvee_{s_j \in F} 0^{|Q|-1-j}10^j$
 - 各マスクテーブルの定義:
 - $B_n[i, \sigma] = \bigvee_{(s_i, \sigma, s_j) \in \Delta} 0^{|Q|-1-j}10^j$
 - $E_n[i] = \bigvee_{s_j \in E(i)} 0^{|Q|-1-j}10^j$ (ここで、 $E(i)$ は状態 s_i の ϵ -closure)
 - $E_d[D] = \bigvee_{i, i=0 \text{ OR } D \& 0^{L-i-1}10^i \neq 0^L} E_n[i]$
 - $B[\sigma] = \bigvee_{i \in 0 \dots m} B_n[i, \sigma]$



擬似コード

```

BuildEps ( $N = (Q_n, \Sigma, I_n, F_n, B_n, E_n)$ )
1  for  $\sigma \in \Sigma$  do
2       $B[\sigma] \leftarrow 0^L$ ;
3      for  $i \in 0 \dots L-1$  do  $B[\sigma] \leftarrow B[\sigma] \mid B_n[i, \sigma]$ ;
4  end of for
5   $E_d[0] \leftarrow E_n[0]$ ;
6  for  $i \in 0 \dots L-1$  do
7      for  $j \in 0 \dots 2^i - 1$  do
8           $E_d[2^i + j] \leftarrow E_n[i] \mid E_d[j]$ ;
9      end of for
10 end of for
11 return ( $B, E_d$ );

```

```

BPTompson ( $N = (Q_n, \Sigma, I_n, F_n, B_n, E_n), T = t_1 t_2 \dots t_n$ )
1  Preprocessing:
2      ( $B, E_d$ )  $\leftarrow$  BuildEps( $N$ );
3  Searching:
4       $D \leftarrow E_d[I_n]$ ;      /* 初期状態 */
5      for  $pos \in 1 \dots n$  do
6          if  $D \ \& \ F_n \neq 0^L$  then report an occurrence ending at  $pos-1$ ;
7           $D \leftarrow E_d[(D \ll 1) \ \& \ B[t_{pos}]]$ ;
8      end of for

```




Bit-parallel Glushkov

G. Navarro and M. Raffinot. Fast regular expression search. *In Proc. of WAE99*, LNCS1668, 199-213, 1999.

- Glushkov NFAをbit-parallelでシミュレートする
 - Glushkov NFAでは、任意の状態について、そこに入る状態遷移のラベルがすべて等しいことに着目
→ Shift-And法的なbit-parallelはできないが、状態遷移は $T_d[D] \& B[\sigma]$ で計算可能
 - マスクテーブルは $2^{|Q|}$ によい (BPTompsonの場合は 2^L)
 - 前処理全体で $O(2^m + m|\Sigma|)$ 時間かかる
 - m が十分小さいときは $O(n)$ 時間でテキストを走査できる
 - ほとんどの場合において、BPTompsonより効率が良い
- NFA $GL=(Q=\{s_0, \dots, s_{|Q|-1}\}, \Sigma, I=s_0, F, \Delta)$ について
 - NFAのマスクビット表現: $Q_n=\{0, \dots, |Q|-1\}$, $I_n=0^{|Q|-1}1$, $F_n=\bigvee_{s_j \in F} 0^{|Q|-1-j}10^j$
 - 各マスクテーブルの定義:
 - $B_n[i, \sigma] = \bigvee_{(s_i, \sigma, s_j) \in \Delta} 0^{|Q|-1-j}10^j$
 - $B[\sigma] = \bigvee_{i \in 0 \dots m} B_n[i, \sigma]$
 - $T_d[D] = \bigvee_{(i, \sigma), D \& 0^{m-i}10^i \neq 0^{m+1}, \sigma \in \Sigma} B_n[i, \sigma]$



擬似コード

BuildTran ($N = (Q_n, \Sigma, I_n, F_n, B_n, E_n)$)

```

1  for  $i \in 0 \dots m$  do  $A[i] \leftarrow 0^{m+1}$ ;
2  for  $\sigma \in \Sigma$  do  $B[\sigma] \leftarrow 0^{m+1}$ ;
3  for  $i \in 0 \dots m, \sigma \in \Sigma$  do
4       $A[i] \leftarrow A[i] \mid B_n[l, \sigma]$ ;
5       $B[\sigma] \leftarrow B[\sigma] \mid B_n[i, \sigma]$ ;
6  end of for
7   $T_d[0] \leftarrow 0^{m+1}$ ;
8  for  $i \in 0 \dots m$  do
9      for  $j \in 0 \dots 2^i - 1$  do
10          $T_d[2^i + j] \leftarrow A[i] \mid T_d[j]$ ;
11     end of for
12 end of for
13 return  $(B, E_d)$ ;

```

BPGlushkov ($N = (Q_n, \Sigma, I_n, F_n, B_n, E_n), T = t_1 t_2 \dots t_n$)

```

1  Preprocessing:
2      for  $\sigma \in \Sigma$  do  $B_n[0, \sigma] \leftarrow B_n[0, \sigma] \mid 0^{m+1}$ ;      /* initial self-loop */
3       $(B, E_d) \leftarrow \mathbf{BuildTran}(N)$ ;
4  Searching:
5       $D \leftarrow 0^{m+1}$ ;      /* 初期状態 */
6      for  $pos \in 1 \dots n$  do
7          if  $D \ \& \ F_n \neq 0^{m+1}$  then report an occurrence ending at  $pos-1$ ;
8           $D \leftarrow T_d[D] \ \& \ B[t_{pos}]$ ;
9      end of for

```



その他のトピックス

■ 拡張正規表現について

– 連結、選択、繰り返しに加えて、積 (intersection) と否定 (complementation) の演算を加えたもの。

■ $\neg(\text{UNIX}) \wedge (\text{UNI}(\cdot)^* \mid (\cdot)^*\text{NIX})$

– 俗に言う (POSIX 定義の) 拡張正規表現とは意味が異なる

■ H. Yamamoto, An Automata-based Recognition Algorithm for Semi-extended Regular Expressions, *Proc. MFCS2000*, LNCS1893, 699–708, 2000.

■ O. Kupferman and S. Zuhovitzky, An Improved Algorithm for the Membership Problem for Extended Regular Expressions, *Proc. MFCS2002*, LNCS2420, 446–458, 2002.

■ 高速化に関する研究

– BNDMを用いたFiltration手法 + 検査

■ G. Navarro and M. Raffinot, New Techniques for Regular Expression Searching, *Algorithmica*, 41(2): 89–116, 2004.

※ この論文には、 $O(m2^m)$ bits のマスクテーブルでGlushkov NFAをシミュレートする方法も併せて記されている



第5回 まとめ

- 正規表現：有限オートマトンと言語を定義する能力が等しい
- 正規表現の照合処理のながれ
 - 構文木 (parse tree) を構築してからNFAへ変換し、NFAをシミュレートしてテキストを走査
 - Filtration+複数パターン照合+検査+NFAシミュレート
- NFAの構築方法
 - ThompsonのNFA
 - 状態数 $< 2m$ 、状態遷移の数 $< 4m \rightarrow O(m)$
 - ϵ 遷移を含む
 - ϵ 遷移以外の遷移は必ず i 番目から $i+1$ 番目に遷移する
 - GlushkovのNFA
 - 状態数はちょうど $m+1$ 個、状態遷移の数は $O(m^2)$
 - ϵ 遷移を含まない
 - 任意のノードについて、そのノードに入ってくる遷移のラベルはすべて等しい
- NFAのシミュレーション手法
 - Thompsonの手法 $\rightarrow O(mn)$ 時間
 - DFAへ変換 $\rightarrow O(n)$ 時間で走査、ただし $O(2^m)$ 時間・領域の前処理
 - Bit-parallel手法による高速化：Bit-parallel Thompson、Bit-parallel Glushkov
- 次回のテーマ
 - 圧縮テキスト上のパターン照合：喜田の研究紹介（本分野における90年代のトレンド！）



付録

- 第一回目の講義で説明していなかった用語の定義について
 - Σ^* の部分集合を**形式言語** (formal language) または単に**言語**という
 - 言語 $L_1, L_2 \in \Sigma^*$ に対して、集合 $\{xy \mid x \in L_1 \text{ かつ } y \in L_2\}$ を L_1 と L_2 の**積** (product) といい、 $L_1 \cdot L_2$ または単に $L_1 L_2$ と書く
 - 言語 $L \subseteq \Sigma^*$ に対して、 $L^0 = \{\varepsilon\}, L^n = L^{n-1} \cdot L \quad (n \geq 1)$ とする。また $L^* = \bigcup_{n=0 \dots \infty} L^n$ として、これを L の**閉包** (closure) という。また $L^+ = \bigcup_{n=1 \dots \infty} L^n$ とする
- 後方参照について
 - 講義の中で、「後方参照については、探してみたが、きちんとした定義がみあたらない」と述べたが、次の文献に記述があることが判明した
 - Handbook of Theoretical Computer Science, Volume A: Algorithms and Complexity, The MIT Press, Elsevier, 1990.
 - (邦訳) コンピュータ基礎理論ハンドブック I : アルゴリズムと複雑さ, 丸善, 1994.
 - 第5章2.3節 後退参照付き正規表現、および6.1節 後退参照付き正規表現の照合問題
 - これによると、後方参照という考え自体は1964年に既に出現していたらしい
 - 正規表現どころか、文脈自由文法の枠組みをも超えている
 - その照合問題は、NP完全 (NP-complete) であることが証明されている